# Unsupervised M·L

Till now we have gone through all supervised learning techniques.

In case of unsupervised ML techniques we have mainly ~~some~~ techniques:

- Clustering techniques

- kNN (K-Nearest neighbours)

- Anamaly detection

- PCA (Principle Component Analysis)

- Neural Network

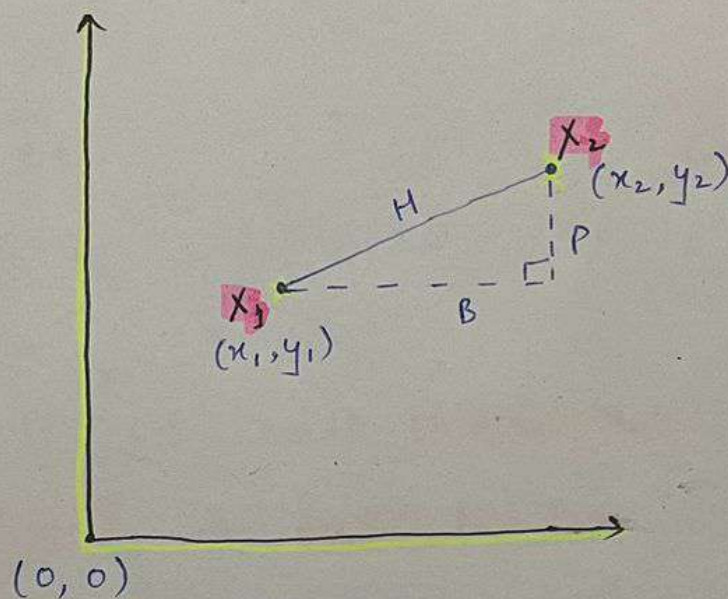- Independent Component Analysis

- Apriari Algorithm.

## Clustering :

It's finding of structure or pattern in collection of uncategarized data, then finding clusters (group) if it exist in the data. we can choose no.s of clusters we want to group our data.

## Deciding datapoint for cluster to belong:

To decide which data point will belong to which cluster we use approach of ==Euclidian distance measure==.

Euclidian distance measure:



distance b/w $X_1$ and $X_2$ can easily be calculated using pythagaras theorem:

$$D_{(x_1, x_2)} = H = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

we also need to decide no.s of clusters bcz we just can't choose ==random no.== as it will ==impact accuracy==.

we use ==Elbow method== to decide no.s of clusters we can choose.

Let's consider below example to have bit idea
of clustering and distance calculation.

| S.no | Height | weight | cluster | BWI |
|------|--------|--------|---------|-----|
| 1 | 170 | 56 | $C_1$ | - |
| 2 | 180 | 63 | $C_1$ | - |
| 3 | 165 | 52 | | - |
| 4 | 176 | 66 | | - |
| 5 | 185 | 78 | $C_2$ | - |
| 6 | 182 | 80 | | - |

$C_1 \rightarrow$ (points to row 2)

$C_2 \rightarrow$ (points to row 5)

Initially we will choose any of point
as centroid of our cluster and then
will compare for which point will fall
in which cluster.

Let's we have consider two clusters, so
we will have two centroid as data of
S.no 2 and 5.

$C_1 \rightarrow$ S.no 2 , $C_2 \rightarrow$ S.no 5

Now using encludian distance approach we
will decide for rest points cluster either
($C_1$ or $C_2$) they will lie.

175

let's consider $C_1$ & $C_2$ as :

$$C_1 \rightarrow (180, 63) \mid C_2 \rightarrow (185, 78)$$

let's consider any of the point from the table say point 1.

Then we will find distance of Point 1 $(P_1)$ from both points $C_1$ & $C_2$

$$d(P_1, C_1) = \sqrt{(170-180)^2 + (56-63)^2}$$

$$P_1 \rightarrow (170, 56) \Big| = 12.2$$

$$d(P_1, C_2) = \sqrt{(170-185)^2 + (56-78)^2}$$

$$= 19.2$$

Now, distance $\boxed{d(P_1, C_1) < d(P_1, C_2)}$

Thus point $P_1$ will lies in the cluster of Centroid $C_1$.

After it lies in the cluster of centroid $C_1$ we will then update the point.

New centroind points of $C_1$ updated :

$$= \left( \frac{170+180}{2}, \frac{56+63}{2} \right)$$
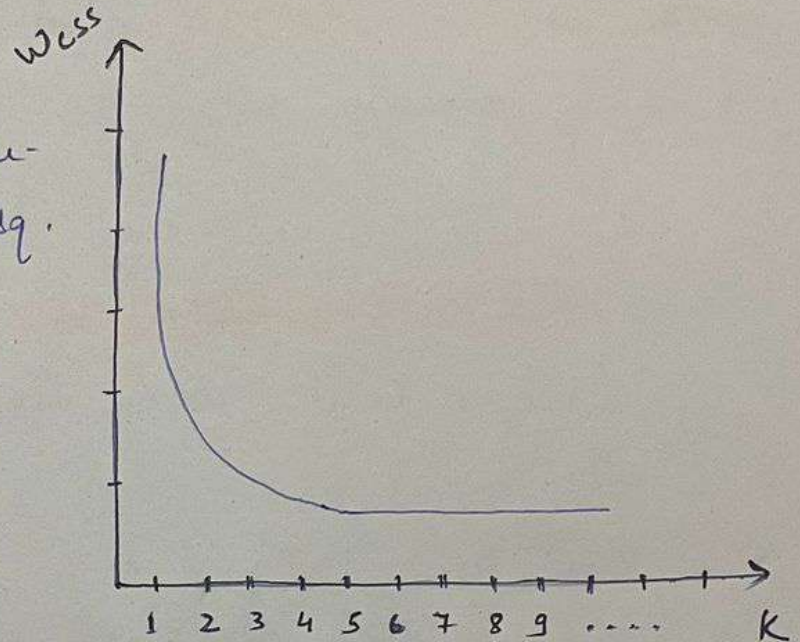
$$= \rightarrow (175, 59.5)$$

176

## ELBOW method:

we use Elbow method to decide no.s of clusters to have best prediction/model.

Wcss: With in clu- -ster sum of sq.

k; No.s of cluster:

$$Wcss = \sum_{j=1}^{n} d(c_j, x_j)$$



$c_j$: Centeroid position / $x_j$: considered datapoint position

Now, Elbow method states that till we have 5 no.s of clusters sum of wcss is huge and after that either we increase no.s of cluster there is almost no diff in wcss.

less wcss will be, more accurate model will be.
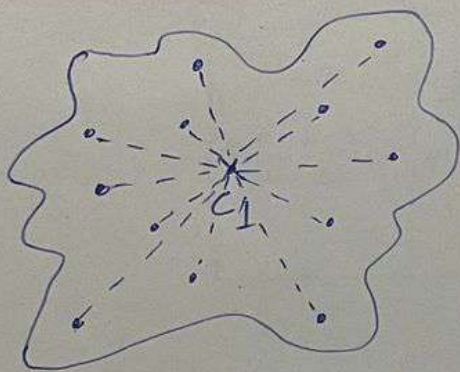
we have two types of clusters of which we find wcss

177

- Intra cluster.
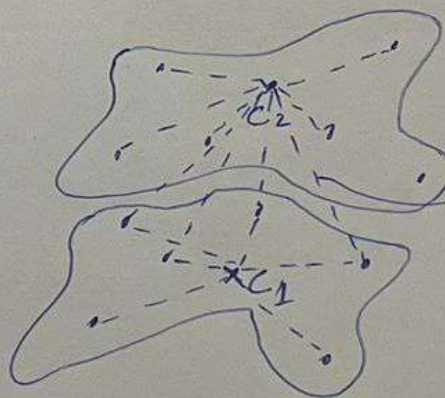- Inter cluster.

## Intra cluster:

when all datapoints lies within the same cluster. we call it intra cluster.

## Inter cluster:

when we have more than one cluster and datapoints lies in each cluster. It's Inter cluster.



Inter cluster calculation.

Intra$^1$ cluster                    Inter$^2$ cluster.

As we know   $wcss = \sum\limits_{j=1}^{n} d(c_j, x_i)^2$

To find wcss square of sum is involved. So more will be distance b/w centroid and datapoint much more will be sq. sum distance.

$$\boxed{wcss_1 >> wcss_2}$$

Thus increasing no.s of cluster decrease wcss sum. But beyond $k = 5$ wcss almost remains same.

## Validation of no.s of clusters (k)

- once we made clusters we need to validate for it scare/accuracy using following methods.

- Dunn Index
- Silhouette scare.

## Dunn Index:

$$\frac{max\ (dis\ (x_i,\ x_j))}{max\ (dis\ (y_i - y_j))}$$

— Inter cluster

— Intra cluster.

## Silhouette scare:

$$\frac{b_j - a_j}{max\ (b_i - a_j)}$$

Inter cluster          Intra cluster.

Silhouette scare lies $b/w$ $-1$ to $1$